

RNA-Seq workflow

Pipeline development

Andrew Ndhlovu

Project SeaStore
Department of Botany and Zoology
Stellenbosch University
<https://github.com/PiscatorX/RNA-Seq-devs>

July 7, 2023



Stellenbosch

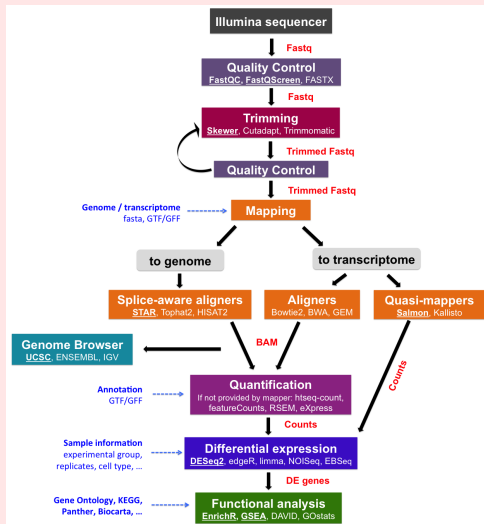
UNIVERSITY
IYUNIVESITHI
UNIVERSITEIT



**PROJECT
SEASTORE**

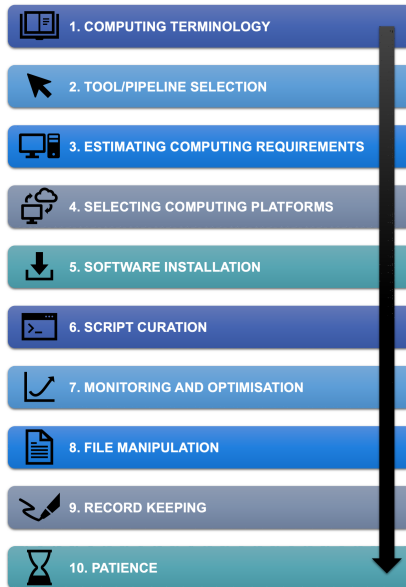
Bioinformatic pipelines and workflow systems

How do we run (and re-run) analyses in way that saves time?



https://biocorecrg.github.io/RNAseq_course_2019/alignment.html

How do we get started



<https://doi.org/10.1371/journal.pcbi.1008645>

Building a bioinformatic pipeline

What is a pipeline

A workflow consisting of a variety of steps (commands) and/or tools that process a given set of inputs to create the desired output files.

Syntax	Paradigm	Interaction	Example	Ease of Development	Ease of Use	Performance
Implicit	Convention	CLI	Snakemake, Nextflow, BigDataScript	★★★★☆	★★★★★	★★★★
Explicit	Convention	CLI	Ruffus, bpipe	★★★★★	★★★★★	★★★★
Explicit	Configuration	CLI	Pegasus	★★★☆☆	★★★★	★★★★★
Explicit	Class	CLI	Queue, Toil	★★★☆☆	★★★★	★★★★★
Implicit	Class	CLI	Luigi	★★★★	★★★★★	★★★★★
Explicit	Configuration	Open Source Server Workbench	Galaxy, Taverna	★★★★	★★★★★	★★★★
Explicit	Configuration	Commercial Cloud Workbench	DNAnexus, SevenBridges	★★★☆☆	★★★★★	★★★★★
Explicit	Configuration	Open Source Cloud API	Arvados, Agave	★★★★	★★★★★	★★★★★

<https://doi.org/10.1093/bib/bbw020>

Most used advanced work flow management systems

nextflow



snakemake
A framework for reproducible data analysis



COMMON
WORKFLOW
LANGUAGE



Criteria	Snakemake	cwltool	Toil	Nextflow
Installation method used	conda	Python pip	Python pip	conda
Number of RiboViz processing steps implemented	14	3	3	5
Time (person-days) to implement steps	1	1	1 ¹	2
Iteration over multiple samples	Yes	See note ²	See note ²	Yes
Error recovery strategies	Yes ("keep going" parameter)	See note ²	See note ²	Yes (step-specific error strategies)
Aggregation of sample-specific results	Yes	See note ²	See note ²	Yes
Conditional execution of steps	Yes (step-specific Python conditions)	No	No ³	Yes (step-specific "when" clauses)
Step-specific log files	Yes (must be captured explicitly by each step)	See note ²	See note ²	Yes (automatically captured)
YAML configuration files	Yes	See note ²	See note ²	Yes
"dry run" option	Yes	See note ²	See note ²	No (but configuration validation can be implemented)
Output a rerunnable bash script to see the commands that were actually executed	No (outputs a summary file with bash commands)	See note ²	See note ²	No (outputs step-specific bash scripts)
Execution within containers, HPC systems, and cloud	Yes	Containers only	Yes	Yes

<https://doi.org/10.1371/journal.pcbi.1008622>

Using a script as pipeline

Save the commands that we run in a file

- Simplest pipeline to implement
- Easiest to learn for beginners
- Usually written in the Unix shell scripting language (Bash)
- Variables and conditional logic can be used for flexibility

Simplicity comes at a cost

- Script pipelines (usually) cannot be continued
- Challenging to implement parallel running of tasks
- Become cumbersome and challenging they become longer
- Still requires some understanding of Linux commands and File system

Introduction to Bash shell scripting

Unix shell is both a command-line interface (CLI) and a scripting language, allowing such repetitive tasks to be done automatically and fast

- Many shells e.g. Bourne shell (sh) and C shell (csh)
- Bash (Bourne Again SHell) is the default shell on Linux
- command are saved in text file i.e the script
- On the command line bash is indicated by the prompt

\$

- Prompt may be preceded by username, machine name and current directory

Learn more

Lots of tutorials online

- The handout from my supervisor is written for people in biology, it is good start:
<https://www.bioinf.wits.ac.za/courses/linux/handout.pdf>
- Highly recommend the scripting tutorial on Tutorials point:
https://www.tutorialspoint.com/unix/shell_scripting.htm